

Scalable MCMC algorithm for the accurate estimation of Exponential Random Graph Models

Byshkin Maksym¹, Alex Stivala², Antonietta Mira¹, Garry Robins², Alessandro Lomi¹

¹Università della Svizzera italiana, Switzerland

²The University of Melbourne, Australia

The increase in the number and size of large networks requires novel efficient methods for their analysis. Exponential Random Graph Models (ERGMs) are exponential family of probability distributions for network structure. The empirical scope of ERGMs is limited by the fact that Maximum Likelihood Estimation (MLE) of model parameters can be obtained only for relatively small networks with a few thousand nodes at most. We propose a new MCMC approach for maximum likelihood estimation of model parameters. In contrast to existing estimation approaches (Bayesian, MCMCMLE and Method of Moments) the approach we propose does not rely on a large number of MCMC simulations to compute equilibrium network configurations. Rather, the Equilibrium Expectations (EE) approach that we propose is based on properties of equilibrium distributions of Markov chains. We implement the EE approach by designing a simple and very efficient algorithm for MLE of ERGM parameters. The suggested estimator is first tested on simulated networks. We compare the results of estimation with results obtained by the Method of Moments and show that the accuracy of estimates obtained via EE is comparable to estimates obtained with the Method of Moments. We demonstrate that the increased speed of the EE algorithm relative to existing methods allows a dramatic increase in the size of networks that can have ERGM parameters estimated by Maximum Likelihood. The empirical value of the EE algorithm is demonstrated also in a study of large biological and social networks.

Maximum Likelihood parameter estimation

$$E_{\pi}(z_A(x)) = z_A(x_{obs})$$

x_{obs} is the observed data

$z_A(x)$ is a set of A functions from the space of states x that measure the properties of the model that are theoretically or empirically relevant (statistics)

$E_{\Theta}(z_A(x))$ are expected statistics that may be obtained by Metropolis-Hastings algorithm if the simulation time is larger than the burn-in time

$$E_{\Theta}(z_A(x)) = \sum_x z_A(x) \pi(x, \Theta)$$

$\pi(x)$ is the probability of state x

$$\pi(x) = \frac{1}{k} \exp\left(-\sum_A \Theta_A z_A(x)\right)$$

Θ_A are model parameters to be estimated

k is a normalizing constant that cannot be computed directly due to complexity

$$k = \sum_x \exp\left(-\sum_A \Theta_A z_A(x)\right)$$

Metropolis-Hastings algorithm

$q(x \rightarrow x')$ proposal

$$\alpha(x \rightarrow x', \Theta) = \min\left\{1, \frac{q(x' \rightarrow x) \alpha(x', \Theta)}{q(x \rightarrow x') \alpha(x, \Theta)}\right\}$$
 acceptance probability

$P(x \rightarrow x', \Theta_A) = q(x \rightarrow x') \alpha(x \rightarrow x', \Theta_A)$ transition probability

Basic sampler

Proposal:
A network dyad is chosen uniformly at random and its value is flipped

$$q(x \rightarrow x') = 1/L_{max}$$

IFD sampler

Proposal: $\left\{ \begin{array}{l} \text{fill random empty dyad} \text{ if } L(x) = L_{obs} \\ \text{delete a random tie} \text{ if } L(x) = L_{obs} + 1 \end{array} \right.$

$$\alpha(x \rightarrow x') = \min\left\{1, \exp\left(\sum_{A \in L} \Theta_A \Delta z_A(x) + \Delta z_L V\right)\right\} \quad V = \left(\Theta_L + \log\left(\frac{L_{max} - L_{obs}}{L_{obs} + 1}\right)\right)$$

Θ_L value is found from Equilibrium Expectation

$$P(L(x) = L_{obs} \rightarrow L(x') = L_{obs} + 1) = P(L(x) = L_{obs} + 1 \rightarrow L(x') = L_{obs})$$

Byshkin M, Stivala A, Mira A, Krause R, Robins G, Lomi A, J. Stat. Phys. 165: 740-754 (2016)

Existing approaches for MLE

- 1) Bayesian
 - 2) Geyer-Thompson MCMC MLE
 - 3) Method of Moments (Stochastic approximation)
- Iteratively modify Θ_A
- At many different Θ_A values perform MCMC simulations to draw equilibrium configurations $x(\Theta_A)$
- MLE is computationally expensive

Geyer CJ & Thompson EA Journal of the Royal Statistical Society. Series B (Methodological), 657-699 (1992)

Snijders, Tom AB. Journal of Social Structure, 1-40 (2002)

New approach for MLE

If x_{obs} is drawn from $\pi(x, \Theta_A^*)$ (simulated network) then we can estimate true Θ^* from

$$\text{“Equilibrium Expectation” (EE): } \sum_x P(x_{obs} \rightarrow x', \hat{\Theta}_A) (z_A(x') - z_A(x_{obs})) = 0$$

We could prove that

If x_{obs} is large $\Theta_A^* = \hat{\Theta}_A$ If x_{obs} is not large $\Theta_A^* = E_{x_{obs}}(\hat{\Theta}_A)$ (Large sample size)

And obtain very fast MLE!

$$x = x_{obs} = \text{constant} \Rightarrow \text{no MCMC simulation}$$

How to apply this theory if x_{obs} is not a simulated network?

MCMC Equilibrium Expectation

$$\left. \begin{array}{l} \sum_x P(x \rightarrow x', \Theta) (z_A(x') - z_A(x)) = 0 \\ z_A(x) = z_A(x_{obs}) \end{array} \right\} \Rightarrow E_{\Theta}(z_A(x)) = z_A(x_{obs})$$

Equilibrium Expectation algorithm

- 1: Initialization: $t=0$; $x = x_{obs}$; $dz_A = 0$; $dz_A(t=0) = 0$
- 2: for $k=1$ to m do
- 3: Propose move $x \rightarrow x'$ with probability $q(x \rightarrow x')$
- 4: Calculate Metropolis-Hastings acceptance probability $\alpha(x \rightarrow x')$
- 5: If $\alpha(x \rightarrow x', \theta(t)) > \text{Unif}([0,1])$ then $dz_A = dz_A + z_A(x') - z_A(x)$ and perform this move: $x = x'$
- 6: end for
- 7: Update of parameters $\theta_A(t+1) = \theta_A(t) - K_A \cdot \text{sgn}(dz_A) (dz_A)^2$
- 8: Increment t . Save sequences $dz_A(t) = dz_A$; If $t < M$ then go to step 2

$$\frac{|dz_A(t > tc)|}{SD(dz_A(t > tc))} < 0.1$$

Simulated networks

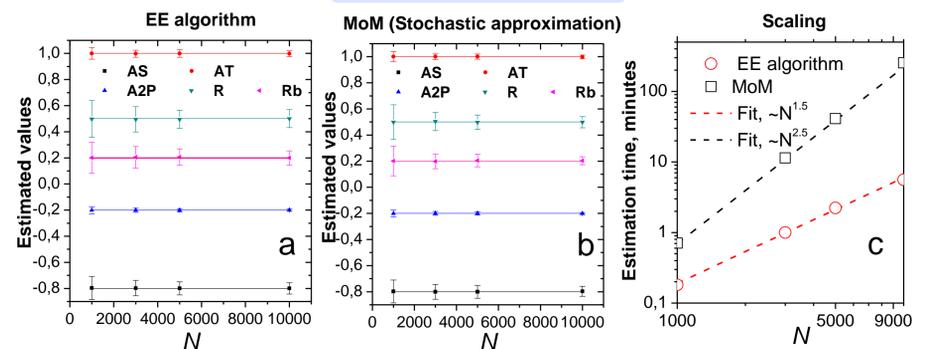


Fig. 1. Estimation of simulated networks by (a) suggested EE algorithm and (b) MoM (Snijders, Tom AB. " Journal of Social Structure (2002): 1-40). The error bars show the confidence intervals, given by doubled values of the standard deviation (over 120 networks). Basic sampler was used. The estimation was performed on Cray XC50 machine available at the Swiss National Supercomputing Centre (CSCS).

AS, AT, A2P, R, Rb are network statistics $z_A(x)$ typical for social networks
Snijders, T.A., Pattison, P.E., Robins, G.L., Handcock, M.S.: New specifications for exponential random graph models. Sociological Methodology 36(1), 99-153 (2006) ; Robins G, Snijders T, Wang P, Handcock M, & Pattison P (2007) Recent developments in exponential random graph (p*) models for social networks. Social Networks 29(2):192-215.

Biological networks

Arabidopsis Interactome Mapping Consortium. Science, 333(6042):601-607, 2011.
C. von Mering, et al. Comparative assessment of large-scale data sets of protein-protein interactions. Nature, 417(6887):399-403, 2002
T. Rolland et al. A proteome-scale map of the human interactome network. Cell, 159(5):1212-1226, 2014.
X.-T. Huang, et al. An integrative C. elegans protein-protein interaction network with reliability assessment based on a probabilistic graphical model. Molecular BioSystems, 12(1):85-92, 2016.
S.-y. Takemura, et al. A visual motion detection circuit suggested by Drosophila connectomics. Nature, 500(7461):175-181, 2013
S. S. Shen-Orr et al. Network motifs in the transcriptional regulation network of Escherichia coli. Nature Genetics, 31(1):64-68, 2002.

Method	Network	Average sample size	N_c	Avg. estim. time (m)	Elapsed time	Estimate (95% C.I.)			
						EE (IFD sampler)	SA (IFD sampler)	SA (basic sampler)	Snowball
EE (IFD sampler)	A. thaliana PPI	2160	20	1.1	01 m 50 s	2.33 (2.32,2.32)	2.32 (2.32,2.32)	2.32 (2.32,2.32)	2.38 (1.72,3.09)
EE (IFD sampler)	Yeast PPI	2617	20	6.6	09 m 07 s	1.28 (1.25,1.31)	1.27 (1.23,1.32)	1.27 (1.23,1.32)	0.00 (-0.07,0.07)
EE (IFD sampler)	Human PPI	4303	20	7.6	10 m 49 s	-14.99 (-15.01,-14.96)	-14.97 (-15.01,-14.96)	-14.97 (-15.01,-14.96)	-14.76 (-16.26,-13.36)
EE (IFD sampler)	C. elegans PPI	5038	20	6.8	09 m 35 s	-7.14 (-7.26,-7.08)	-7.12 (-7.26,-7.08)	-7.12 (-7.26,-7.08)	-10.49 (-11.51,-9.49)
EE (IFD sampler)	E. coli regulatory	418	20	0.6	00 m 43 s	-0.05 (-0.05,0.05)	-0.05 (-0.05,0.05)	-0.05 (-0.05,0.05)	-0.57 (-0.61,0.29)
EE (IFD sampler)	Drosophila optic medulla	1781	20	4.3	06 m 22 s	1.86 (1.81,1.91)	1.86 (1.82,1.90)	1.86 (1.82,1.90)	0.85 (0.18,1.07)
SA (IFD sampler)	A. thaliana PPI	2160	20	9.2	01 h 34 m 02 s	-7.78 (-7.81,-7.71)	-7.78 (-7.81,-7.71)	-7.78 (-7.81,-7.71)	-6.57 (-13.42,-4.88)
SA (IFD sampler)	Human PPI	4303	20	49.2	2 h 46 m 54 s	1.32 (1.29,1.35)	1.32 (1.29,1.35)	1.32 (1.29,1.35)	1.29 (0.82,1.16)
SA (IFD sampler)	Yeast PPI	2617	20	45.6	2 h 02 m 38 s	1.37 (1.36,1.38)	1.37 (1.36,1.38)	1.37 (1.36,1.38)	0.35 (0.02,0.08)
SA (IFD sampler)	C. elegans PPI	5038	20	766.5	25 h 07 m 44 s	-11.03 (-11.08,-10.98)	-11.03 (-11.08,-10.98)	-11.03 (-11.08,-10.98)	-8.82 (0.02,0.07)
SA (IFD sampler)	E. coli regulatory	418	20	0.0	0 h 00 m 06 s	-11.77 (-11.82,-11.73)	-11.77 (-11.82,-11.73)	-11.77 (-11.82,-11.73)	-9.04 (-13.21,-7.31)
SA (IFD sampler)	Drosophila optic medulla	1781	20	824.6	72 h 40 m 00 s	1.04 (1.01,1.07)	1.04 (1.01,1.07)	1.04 (1.01,1.07)	1.18 (0.82,1.16)
SA (basic sampler)	A. thaliana PPI	2160	0	—	(time limit)	1.59 (1.58,1.61)	1.59 (1.57,1.61)	1.59 (1.57,1.61)	0.35 (0.02,0.08)
SA (basic sampler)	Yeast PPI	2617	0	—	(time limit)	-11.03 (-11.08,-10.98)	-11.03 (-11.08,-10.98)	-11.03 (-11.08,-10.98)	-8.82 (0.02,0.07)
SA (basic sampler)	Human PPI	4303	0	—	(time limit)	0.45 (0.42,0.48)	0.44 (0.41,0.47)	0.44 (0.41,0.47)	0.44 (0.18,0.69)
SA (basic sampler)	C. elegans PPI	5038	3	204.5	7 h 40 m 20 s	0.78 (0.64,0.93)	0.79 (0.66,0.92)	0.79 (0.66,0.92)	0.79 (0.61,0.96)
SA (basic sampler)	E. coli regulatory	418	20	1.1	0 h 04 m 06 s	-6.55 (-6.55,-6.55)	-6.53 (-6.53,-6.53)	-6.53 (-6.53,-6.53)	-6.53 (-6.53,-6.53)
SA (basic sampler)	Drosophila optic medulla	1781	0	—	(time limit)	0.23 (0.17,0.30)	0.24 (0.18,0.30)	0.24 (0.18,0.30)	1.17 (0.82,1.58)
Snowball sampling	A. thaliana PPI	490.6	19	26.3	21 h 08 m 24 s	—	—	—	—
Snowball sampling	Yeast PPI	264.8	19	30.2	3 h 40 m 34 s	—	—	—	—
Snowball sampling	Human PPI	822.5	18	47.0	3 h 50 m 27 s	—	—	—	—
Snowball sampling	C. elegans PPI	496.4	16	270.7	40 h 00 m 33 s	—	—	—	—
Snowball sampling	E. coli regulatory	649.7	15	118.0	7 h 22 m 48 s	—	—	—	—
Snowball sampling	Drosophila optic medulla	—	—	—	—	—	—	—	—

A. Stivala, et al. Snowball sampling for estimating exponential random graph models for large networks. Social Networks 47, 167-188 (2016)

Large social network

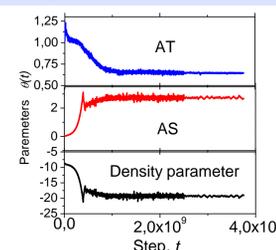


Fig. 2. MLE of ERGM parameters for Livemocha networks with 104103 nodes and 2193083 ties

More presentations

- Complex Networks, March 20-24, 2017, Dubrovnik, Croatia
- Sunbelt INSNA Conference, May 30-June 4, 2017, Beijing, China
- Cambridge Networks Day, 13th June 2017, Cambridge, UK
- PASC17, June 26 - 28, 2017, Lugano, Switzerland
- International Conference on Monte Carlo Methods and Applications, July 3-7, 2017, Montreal, Canada
- International Conference on Computational Social Science, July 10-13, 2017, Cologne, Germany
- Third European Conference on Social Networks, September 26 to 29, 2017, Mainz, Germany

Zafarani and H. Liu. Social computing data repository at ASU, 2009
Livemocha network datasets, konekt.uni-koblenz.de